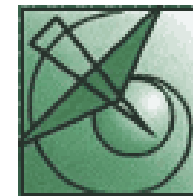




# Odborná terminologie knihovnictví a informační vědy očima uživatelů databáze TDKIV

## Předběžné výsledky projektu

Helena Kučerová  
VOŠIS Praha





# Cíl projektu:

Využít metody **kvantitativní analýzy** k objevení nových poznatků a námětů, skrytých v uchovávaných **uživatelských datech** o vyhledávání v TDKIV.

- Zadavatel:  
Redakční rada TDKIV
- Řešitel:  
VOŠIS Praha – 6členný studentský tým
- Termín: zimní semestr 2011/2012





# TDKIV v číslech:

rok vzniku: **2003**

počet (preferovaných) termínů: **3 000**

počet ekvivalentů (nepreferovaných termínů):  
**4 000**

počet dotazů / rok: **53 000**

počet vyhledávaných výrazů / rok: **12 000**

počet uživatelů (UIP adres): **500**

průměrná doba strávená v databázi: **1 – 2 minuty**



# Struktura terminologického hesla

- **termín**
- anglický ekvivalent
- **ekvivalent** (nepreferovaný termín)
- **příbuzný termín**
- výklad termínu (autorský / normativní)
- zdroj výkladu
- poznámka
- třídník (věcná kategorie)
- autor / konzultant / lektor / redaktor



# Ukázka terminologického hesla

<b>Termín</b>	• <a href="#">data mining</a>
<b>Poznámka</b>	ciz., neol., český ekvivalent není ustálený
<b>Ekvivalent</b>	<b>dolování dat</b> (neol.) text mining vytěžování informací vytěžování médií vytěžování obsahu databáze web mining (ciz.)
<b>Termín anglicky</b>	data mining
<b>Výklad termínu</b>	Technologie vyhledávání, modelování a prezentace předem neznámých informací, příp. znalostí a vztahů mezi daty v rozsáhlých databázích a datových skladech. Analýzy se odvozují přímo z obsahu <b>dat</b> , nikoliv na základě hypotéz či dotazů uživatele. Využívají se techniky umělé inteligence (neuronové sítě, rozpoznávání, samoučící se algoritmy), jež mohou být kombinovány s technikami statistického a matematického modelování (klasifikační pravidla nebo stromy, regrese, shluková analýza) a s nástroji OLAP (on-line analytické zpracování). [KUPKA-2002:57-60] [POKORNÝ-1999:96-99]
<b>Zdroj výkladu</b>	<a href="#">KUPKA-2002:57-60</a> <a href="#">POKORNÝ-1999:96-99</a>
<b>Příbuzný termín</b>	• <a href="#">datový sklad</a> • <a href="#">extrakce informací</a> • <a href="#">knowledge discovery in databases</a>
<b>Viz též</b>	• <a href="#">záznam v bázi KSL</a>
<b>Autor hesla</b>	Kučerová, Helena
<b>Lektor hesla</b>	Kimlička, Štefan
<b>Redaktor hesla</b>	Schwarz, Josef Burgetová, Jarmila
<b>Systém. číslo</b>	000000088



# Možnosti vyhledávání v TDKIV

Způsob vyhledávání	Použitý index	Zdrojová báze
vyhledávání SEARCH	z více polí FREE TEXT	KTD
listování SCAN		
navigace Příbuzný termín	z jednoho pole	KTDP Pracovní báze



# Vyhledávání v TDKIV

Databáze Národní knihovny ČR Aktuální báze: KTD

	Konec	Přihlášení	Databáze	Díličí báze	Nastavení	Otázky	Nápověda
	Vyhledávání / Rejstříky	Výsledky dotazu	Předchozí dotazy	Schránka	Historie		

Základní vyhledávání | Z víceází | Vyhledávání CCL | Nový záznam

## KTD - Česká terminologická databáze knihovnictví a informační vědy (TDKIV)

[Vstup do pracovní báze](#)

[Základní vyhledávání](#)

[Informace o excerptci z norem](#)

[Informace o bázi](#)

SEARCH

Vyberte údaj pro vyhledávání:	Free-text	<input type="text"/>	<input type="button" value="OK"/>
Zadejte slovo nebo slovní spojení:	Free-text	<input type="text"/>	<input type="button" value="Vyčistit formulář"/>
Blízkost slov?	Termín / ekvivalent	<input type="text"/>	
	Třídník	<input type="text"/>	
	Systémové číslo	<input type="text"/>	

[Prohlížení rejstříků](#)

Vyberte rejstřík k prohlížení:	Termín / Ekvivalent	<input type="button" value="OK"/>
Zadejte heslo k otevření rejstříku:	<input type="text"/>	<input type="button" value="Vyčistit formulář"/>

[Nápověda k vyhledávání:](#)

- V dotazu nezáleží na tom, zda použijete velká nebo malá písmena.
- Vyhledávání "Free-text" hledá příslušné slovo/fetězec v polích *Termín/Ekvivalent*, *Anglický ekvivalent*, *Výklad a Poznámka*. K vyhledávání podle "Třídníků" můžete využít [pracovní tabulku](#). K vyhledávání podle "Statusu záznamu" použijte [tabulku statusů](#).





# Listování v TDKIV

Databáze Národní knihovny ČR Aktuální báze: KTD (prac.)

<b>NK</b>	Konec	Přihlášení	Databáze	Dílčí báze	Nastavení	Otázky	Nápověda
	Vyhledávání / Rejstříky	Výsledky dotazu	Předchozí dotazy	Schránka	Historie	🇬🇧	

Základní vyhledávání | Rozšířené vyhledávání | Vyhledávání CCL | Nový záznam

**KTD - Pracovní verze**  
Česká terminologická databáze knihovnictví a informační vědy  
(TDKIV)

[Návrat do báze KTD \(TDKIV\)](#)

Základní vyhledávání

[Informace o bázi](#)

Vyberte údaj pro vyhledávání:	Free-text	<input type="button" value="OK"/> <input type="button" value="Vyčistit formulář"/>
Zadejte slovo nebo slovní spojení:	<input type="text"/>	
Blízkost slov?	<input checked="" type="radio"/> Ne <input type="radio"/> Ano	

SCAN

Prohlížení rejstříků

Vyberte rejstřík k prohlížení:	Termín / Ekvivalent	<input type="button" value="OK"/> <input type="button" value="Vyčistit formulář"/>
Zadejte heslo k otevření rejstříku:	<input type="text"/>	
Nápověda k vyhledávání:	<ul style="list-style-type: none"><li>Termín / Ekvivalent</li><li>Anglický ekvivalent</li><li>Třídník</li><li>Autor hesla</li><li>Konzultant</li><li>Lektor hesla</li><li>Zdroj výkladu</li></ul>	

- V dotazu nezáleží na tom, zda použijete velká
- Vyhledávání "Free-text" hledá příslušné slovořetec v poli *Termín/Ekvivalent*, *Anglický ekvivalent*, *Výklad a Poznámka*. K vyhledávání podle "Třídniců" můžete využít [pracovní tabulku](#). K vyhledávání podle "Statusu záznamu" použijte [tabulku statusů](#).



# Navigace v TDKIV

<b>Termín</b>	• <a href="#">řízený slovník</a>
<b>Termín anglicky</b>	controlled vocabulary
<b>Výklad termínu</b>	<b>Slovník</b> lexikálních jednotek selekčního jazyka uspořádaný specifickým způsobem (např. zahrnuje vztahy ekvivalence, hierarchie a asociace), který slouží pro indexaci a vyhledávání dokumentů. [KATUŠČÁK, MATTHAEIDESOVÁ, NOVÁKOVÁ-1998:306]
<b>Normativní výklad</b>	Seznam slov nebo frází, které byly schváleny pro indexování. [ČSN ISO 5127-2003]
<b>Zdroj výkladu</b>	<a href="#">KATUŠČÁK, MATTHAEIDESOVÁ, NOVÁKOVÁ-1998:306</a>
<b>Zdroj norm.výkladu</b>	<a href="#">ČSN ISO 5127-2003</a>
<b>Příbuzný termín</b>	• <a href="#">kvalifikovaný Dublin Core</a> • <a href="#">nepreferovaný termín</a> • <a href="#">předmětový heslář</a> • <a href="#">syntetická metoda</a> • <a href="#">tezaurus</a>
<b>Autor hesla</b>	Balíková, Marie
<b>Lektor hesla</b>	Hrazdil, Aleš
<b>Redaktor hesla</b>	Schwarz, Josef Burgetová, Jarmila
<b>System. číslo</b>	000001624



# Výzkumné otázky:

## **Statistický rozbor a vizualizace:**

- četnost hledaných a nalezených / nenalezených výrazů
- četnost způsobů vyhledávání
- čas vyhledávání
- rozdělení IP adres

## **Data mining a business intelligence:**

- příčiny neúspěšných dotazů
- společně hledané termíny



# Metodika:

- korpusová terminografie
- kvantitativní analýzy, pokus o data mining (bibliomining)
- logy z vyhledávání v systému Aleph – data za období **březen 2010 – srpen 2011**  
cca **80.000** záznamů
- MS SQL, dotazování v SQL



2010030106062306	217.197	1	21	KTD	Free-text= sigla
2010030108372010	195.113	6	21	KTD	Free-text= informační management
2010030108413077	147.229	0	21	KTD	Free-text= avízo
2010030108423950	195.113	2	21	KTD	Free-text= ontologie
2010030108425062	195.113	2	21	KTD	Free-text= ontologie
2010030108425941	195.113	6	21	KTD	Free-text= informační management
2010030108431765	195.113	6	21	KTD	Free-text= informační management
2010030108491647	195.113	1	21	KTD	Free-text= certifikační autority
2010030108501060	195.113	1	21	KTD	Free-text= certifikační autority
2010030108510436	195.113	1	21	KTD	Free-text= certifikační autority
2010030108511091	195.113	1	21	KTD	Free-text= rss

```

SELECT a.search_text as hledany_a,
COUNT(a.search_text) as pocet_hledani,
b.search_text as hledany_b
FROM SEARCH_KTD a
JOIN SEARCH_KTD b on a.cas = b.cas
AND a.ip = b.ip and a.datum = b.datum
WHERE a.search_text <> b.search_text
AND a.HITS <> 0 and b.HITS <> 0
GROUP BY a.search_TEXT, b.SEARCH_TEXT,
a.SEARCH_TEXT+a.DATUM+a.CAS,
b.SEARCH_TEXT+b.DATUM+b.CAS
HAVING COUNT(a.SEARCH_TEXT) > 5

```

eriálové publikace  
ntologie  
olandr  
olandr  
orma  
OAN PERIOD  
OAN  
apír  
ýroba papíru  
apírenství  
apírenství  
apírenství  
oan period  
ace  
ourdni  
nihtisk  
nižnično-informačné služby  
konómie model IS-LM  
nižnično-informačné služby

2010030109390139	148.182	0	21	KTD	termin- ekonomie model IS-LM
2010030109594718	158.195	0	21	KTD	Free-text= knižnično-informačné služby
2010030110010636	158.195	25	21	KTD	Free-text= informační služby
2010030110014915	158.195	25	21	KTD	Free-text= informační služby
2010030110015666	158.195	0	21	KTD	Free-text= knižniční služby
2010030110025765	158.195	0	21	KTD	Free-text= knižniční služby
2010030110035919	158.195	0	21	KTD	Free-text= knižniční služby
2010030110110749	010.001	1	21	KTD	Free-text= certifikáty
2010030110355599	195.113	6	21	KTD	Free-text= kognitivní



# Struktura dat

čas zadání dotazu

část IP adresy počítače,  
ze kterého byl dotaz zadán

počet nalezených záznamů

způsob vyhledávání

(21 – základní vyhledávání,  
23 – pokročilé vyhledávání,  
29 – vyhledávání v rejstřících...)

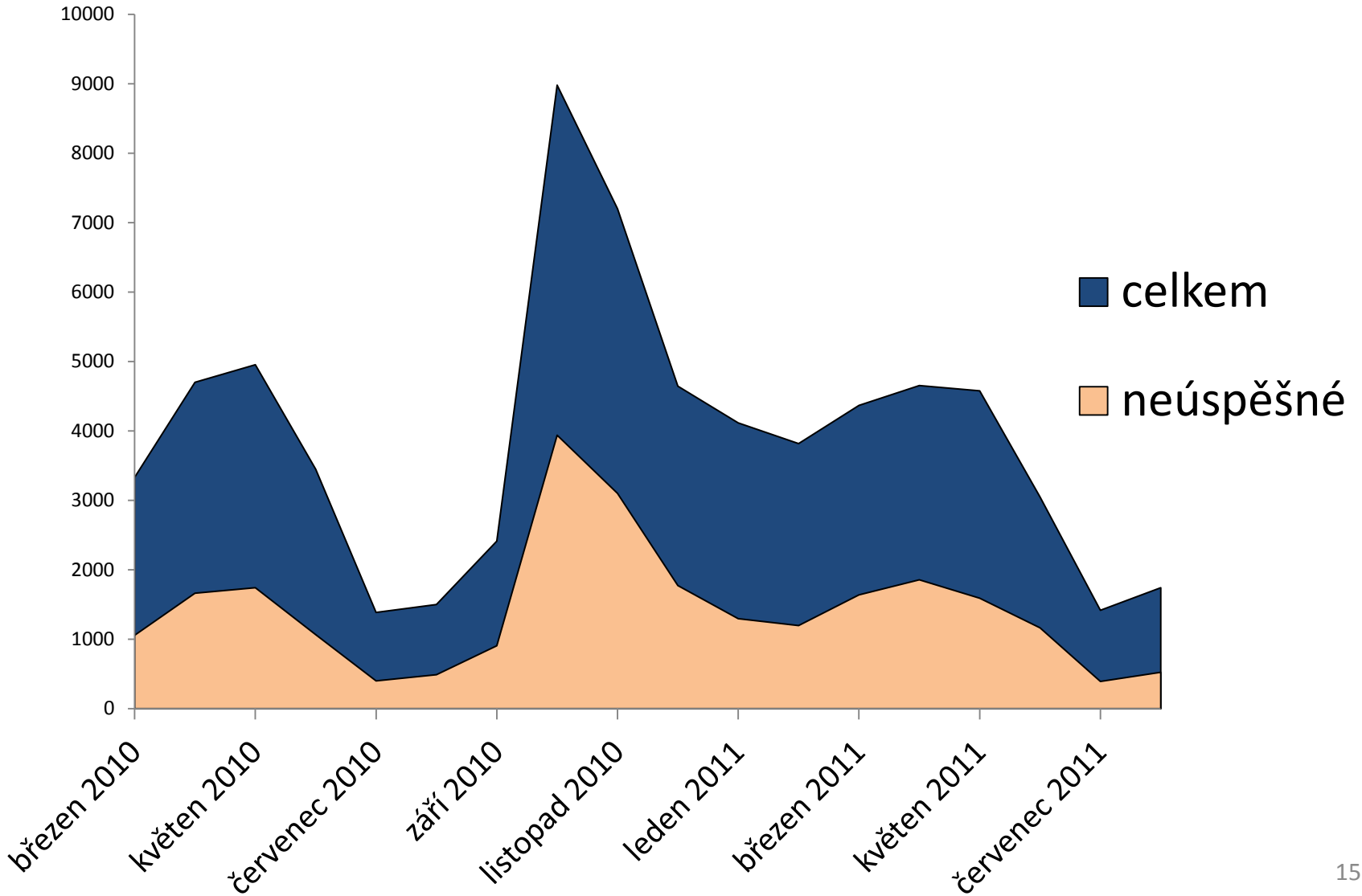
prohledávaná báze  
(KTD, KTDP)

vyhledávaný výraz

TIME_STAMP	IP	HITS	TYP	SEARCH	SEARCH_TEXT
2010030106062306	217.197	1	21	KTD	Free-text= sigla
2010030108372010	195.113	6	21	KTD	Free-text= informační management
2010030108413077	147.229	0	21	KTD	Free-text= avizo
2010030108423950	195.113	2	21	KTD	Free-text= ontologie
2010030108425062	195.113	2	21	KTD	Free-text= ontologie
2010030108425941	195.113	6	21	KTD	Free-text= informační management
2010030108431717	195.113	6	21	KTD	Free-text= informační management
2010030108491647	195.113	1	21	KTD	Free-text= certifikační autority
2010030108501060	195.113	1	21	KTD	Free-text= certifikační autority
2010030108510436	195.113	1	21	KTD	Free-text= certifikační autority
2010030108511091	195.113	1	21	KTD	Free-text= rss
2010030109014741	147.229	9	21	KTD	Free-text= seriálové publikace
2010030109030701	193.084	2	21	KTD	Free-text= ontologie
2010030109093525	193.084	0	21	KTD	Free-text= holandr
2010030109100751	193.084	0	21	KTD	Free-text= holandr
2010030109145152	195.113	6	21	KTD	Free-text= norma
2010030109151268	195.113	1	21	KTD	Free-text= LOAN PERIOD
2010030109170096	195.113	25	21	KTD	Free-text= LOAN
2010030109191401	193.084	73	21	KTD	Free-text= papír
2010030109193111	193.084	2	21	KTD	Free-text= výroba papíru
2010030109195423	193.084	0	21	KTD	Free-text= papírenství
2010030109200318	193.084	0	21	KTD	Free-text= papírenství
2010030109201871	193.084	0	21	KTD	Free-text= papírenství
2010030109202824	195.113	1	21	KTD	Free-text= loan period
2010030109230840	193.084	11	21	KTD	Termin= anotace
2010030109352534	193.084	0	21	KTD	Free-text= fourdni
2010030109400995	193.084	2	21	KTD	Free-text= knihtisk
2010030109574758	158.195	0	21	KTD	Free-text= knižnično-informačné služby
2010030109584346	146.102	0	21	KTD	Free-text= ekonomie model IS-LM
2010030109584714	158.195	0	21	KTD	Free-text= knižnično-informačné služby
2010030109590139	146.102	0	21	KTD	Termin= ekonomie model IS-LM
2010030109594718	158.195	0	21	KTD	Free-text= knižnično-informačné služby
2010030110010636	158.195	25	21	KTD	Free-text= informační služby
2010030110014915	158.195	25	21	KTD	Free-text= informační služby
2010030110015666	158.195	0	21	KTD	Free-text= knižniční služby

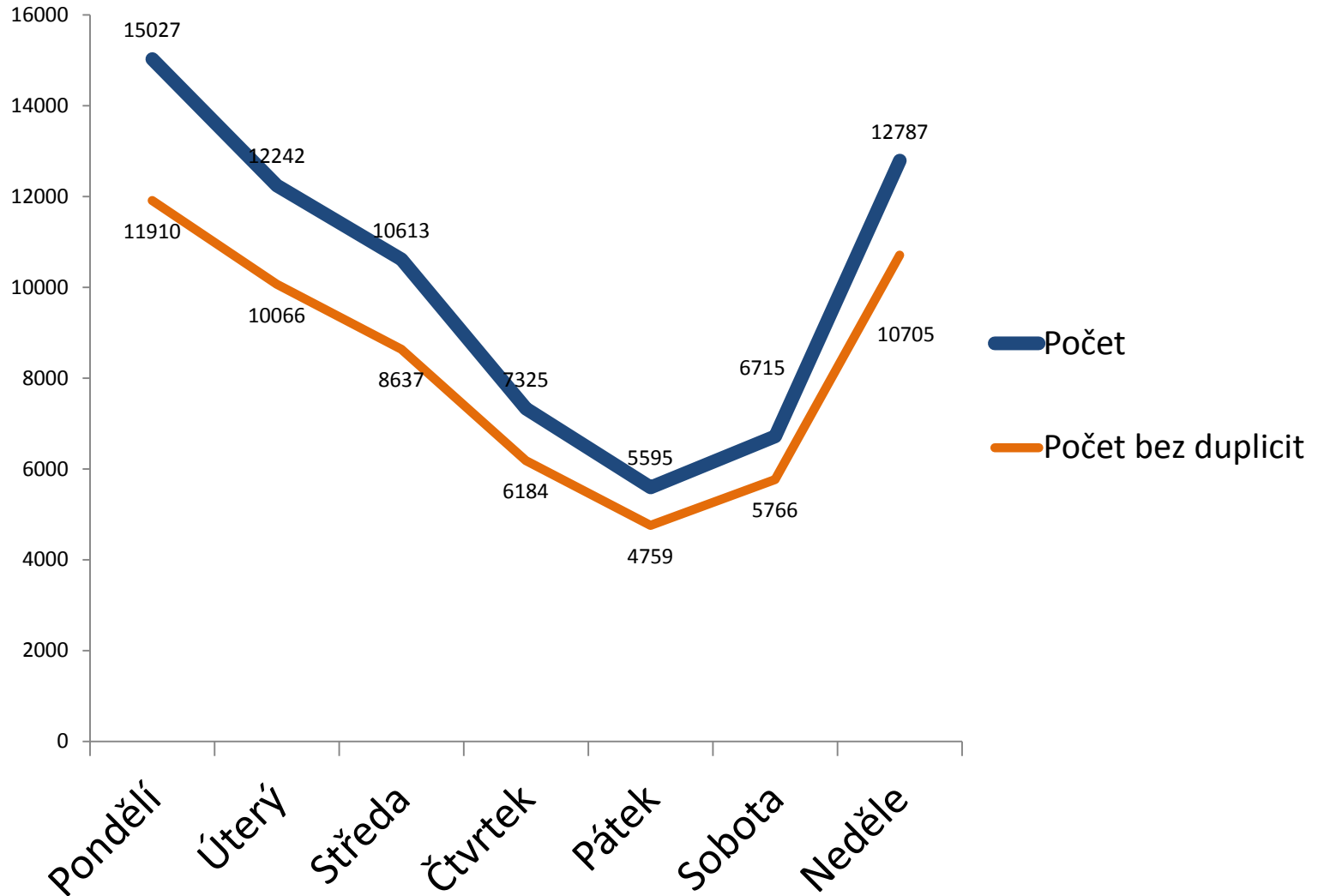


# Vyhledávání podle měsíců





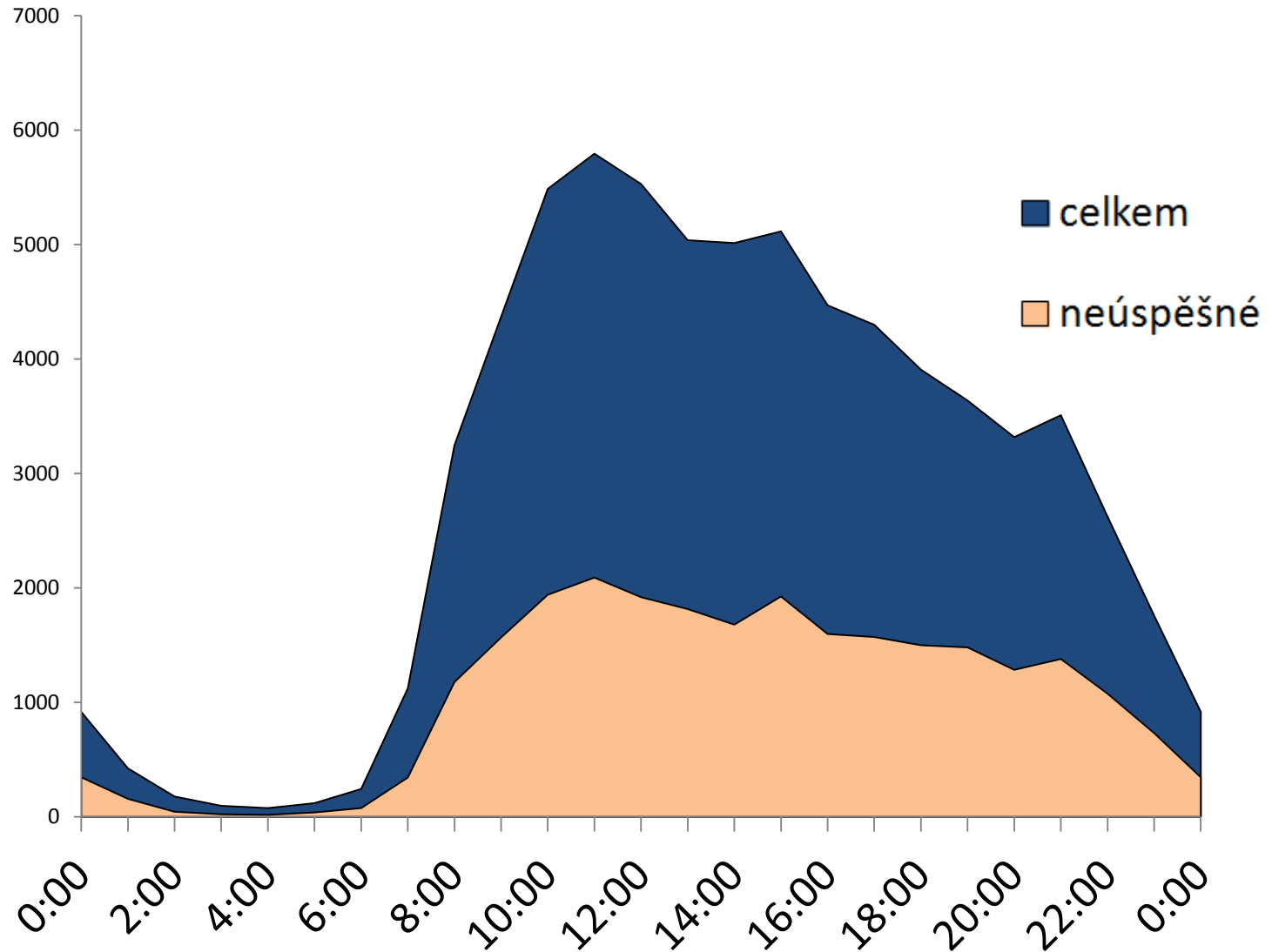
# Vyhledávání podle dní







# Čas vyhledávání





# Rozdělení uživatelů

Vyhledávání v KTD:

42 701 vyhledávání (sessions)

4 264 (cca 10 %) vyhledávání z IP 195.113.xx (PASNET)

	počet vyhledávání		počet IP adres	
SCAN	6 026 80 %	↑ <b>80 %</b>	50 21 %	↓ <b>20 %</b>
SEARCH	28 134 80 %		112 24 %	
SCAN	1 507 20 %	↑ <b>20 %</b>	181 80 %	↓ <b>80 %</b>
SEARCH	7 034 20 %		361 76 %	



# Nejvyhledávanější výrazy

Free-text= akvizice	432
Free-text= Rešerše	389
Free-text= informace	370
Free-text= knihovna	348
Free-text= bibliografie	264
Free-text= dokument	228
Free-text= anotace	225
Free-text= Databáze	200
Free-text= abstrakt	199
Free-text= tezaurus	195
Free-text= Informační věda	194
Free-text= informační gramotnost	180
Termín= informace	177
Free-text= kniha	164
Free-text= informační zdroj	159
Free-text= signatura	155
Free-text= monografie	154
Free-text= Informační služby	150
Free-text= katalog	141
Termín= dokument	133



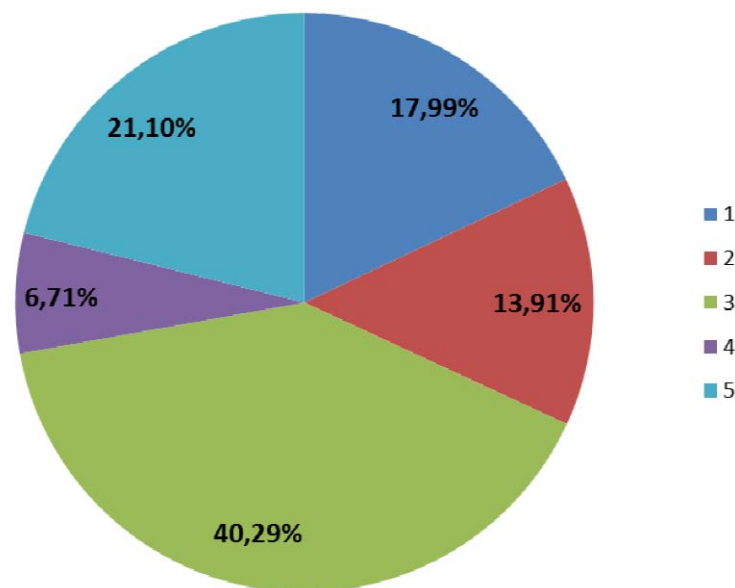
# Úspěšnost vyhledávání

<b>počet vyhledávaných výrazů</b>	<b>17 285</b>
úspěch ( $\geq 1$ hitů)	<b>6 867</b> (cca 40 %)
neúspěch (0 hitů)	<b>10 326</b> (cca 60 %)



# Nenalezené výrazy

Free-text= informační potřeba	104
Free-text= AACR2	59
Free-text= signet	57
Free-text= informačný zdroj	49
Free-text= čtenářský fond	39
Free-text= sociální síť	39
Free-text= informační architektura	38
Free-text= registre	37
Free-text= pagina	37
Termín= Registre	36
Free-text= sfx	33
Termín= bibliografia	32
Free-text= jednotná informační brána	30
Free-text= pozvánka	30
Free-text= taxonomie	30
Free-text= facebook	28
Free-text= elzevirský formát	28
Free-text= inkunábula	27
Free-text= česká národní bibliografie	26
Free-text= repozitář	26



1 překlep

2 doporučeno zařadit do TDKIV

3 věcně nerelevantní k TDKIV

4 nesprávný typ vyhledávání

5 jiný jazyk než čeština



# Výrazy hledané společně

pořádání informací	informační analýza	143
primární dokument	dokument	91
popisný údaj	bibliografický popis	70
referát	anotace	70
věcné zpracování	indexace	66
katalogizátor	katalogizační záznam	55
abstrakt	anotace	54
informační gramotnost	databáze	48
zpracování dokumentů	dokument	42
deskriptor	tezaurus	40
synopse	referát	30
informační analýza	pořádání informací	26
jmenné zpracování	věcné pořádání	24



# Předběžné shrnutí výsledků

**Potvrzeno:** Kvantitativní metody mohou poskytnout jak náměty ke zkvalitnění použitelnosti a přístupnosti databáze, tak i přímé podněty k terminologické práci

- náměty na zařazení nových termínů nebo ekvivalentů do TDKIV
- klastry současně hledaných termínů umožňují nacházet sémantické vztahy
- odhaleny problémy k řešení:
  - 2/3 vyhledávání končí neúspěchem
  - 1/5 z nenalezených výrazů jsou překlepy
  - významný podíl dotazů ve slovenštině
  - dotazy na zkratky

erata  
inforamční průmysl  
aldinky



# A co na to studenti...



*Rozšířili jsme si znalosti programu MS SQL, získali větší praxi s prací s databázemi, získali jsme informace o dataminingu – co to je, jak se používá, naučili jsme se, že díky dataminingu můžeme z databáze, kterou máme k dispozici získat takové informace, o kterých jsme předtím ani nevěděli, že je můžeme zjistit.*

*Co jsou lidé vůbec schopni napsat do vyhledávače Národní knihovny,*

*Zjistil jsem, že se dataminingem už v budoucnosti nechci zabývat, není to obor pro mě.*

*Tento projekt mi změnil názor na to, co všechno se dá s daty dělat. Taky to, že z takového malého množství sloupců (kategorií) se dá vyvodit mnoho informací užitečných, které nejdou na první pohled vidět.*





# Plány na další pokračování výzkumu:

- tvorba vícečetných klastrů
- ontologizace TDKIV  
(syntagmatické → paradigmatické vztahy)
- porovnání vyhledávacích výrazů s termíny v databázi
- analýza vyhledávacích procesů
- „opravdový“ data mining 😊

**Další náměty?**